

Regression

RAD Lab Machine Learning Workshop

Kurt Miller

08/23/07

Adapted from slides by Romain Thibaux

Outline

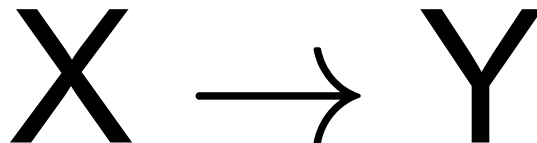
- Ordinary Least Squares Regression
 - Online version
 - Normal equations
 - Probabilistic interpretation
- Overfitting and Regularization
- Overview of additional topics
 - L_1 Regression
 - Quantile Regression
 - Kernel Regression and LWR
 - Spline Regression

Outline

- Ordinary Least Squares Regression
 - Online version
 - Normal equations
 - Probabilistic interpretation
- Overfitting and Regularization
- Overview of additional topics
 - L_1 Regression
 - Quantile Regression
 - Kernel Regression and LWR
 - Spline Regression

Regression vs. Classification:

Classification



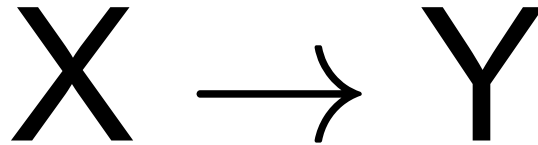
Anything:

- continuous (\mathfrak{R} , \mathfrak{R}^d , ...)
- discrete ($\{0,1\}$, $\{1,\dots,k\}$, ...)
- structured (tree, string, ...)
- ...

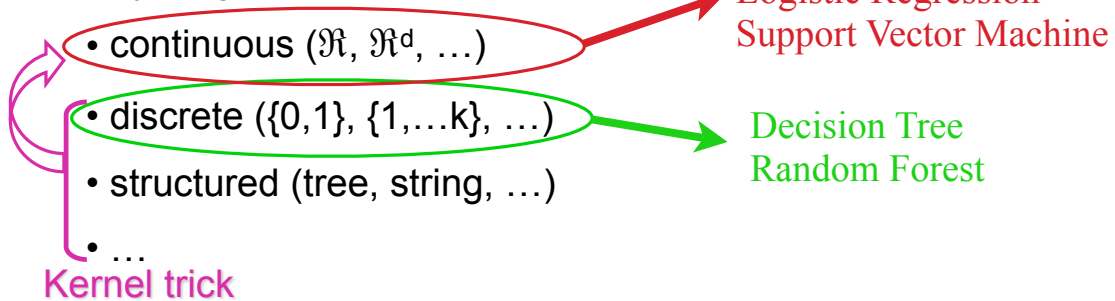
• discrete:

- $\{0,1\}$ *binary*
- $\{1,\dots,k\}$ *multi-class*
- tree, etc. *structured*

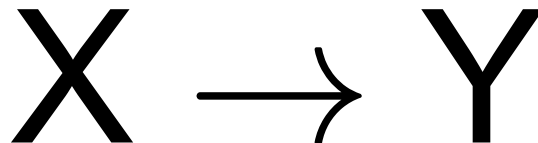
Regression vs. Classification: Classification



Anything:



Regression vs. Classification: Regression



Anything:

- continuous ($\mathcal{R}, \mathcal{R}^d, \dots$)
 - discrete ($\{0,1\}, \{1, \dots, k\}, \dots$)
 - structured (tree, string, ...)
 - ...
- continuous:
– $\mathcal{R}, \mathcal{R}^d$

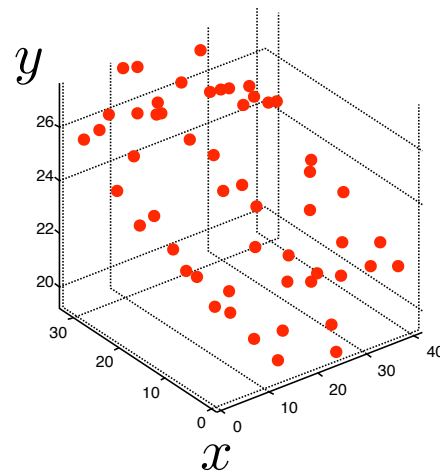
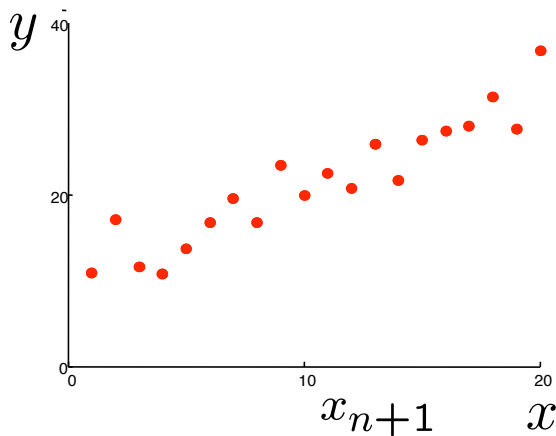
Examples

- Voltage \rightarrow Temperature
- Processes, memory \rightarrow Power consumption
- Protein structure \rightarrow Energy
- Robot arm controls \rightarrow Torque at effector
- Location, industry, past losses \rightarrow Premium

Linear regression

Given examples $(x_i, y_i)_{i=1\dots n}$

Predict y_{n+1} given a new point x_{n+1}

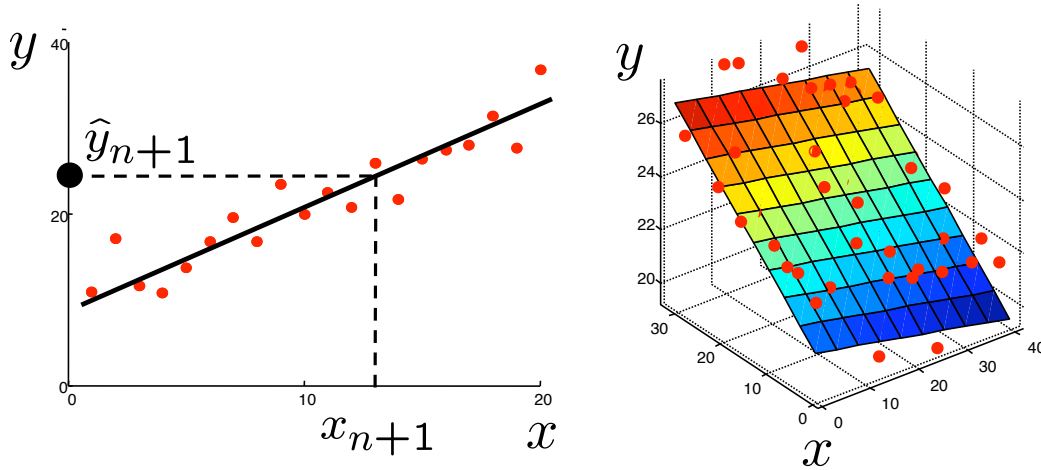


Linear regression

We wish to estimate \hat{y} by a linear function of our data x :

$$\begin{aligned}\hat{y}_{n+1} &= w_0 + w_1 x_{n+1,1} + w_2 x_{n+1,2} \\ &= w^\top x_{n+1}\end{aligned}$$

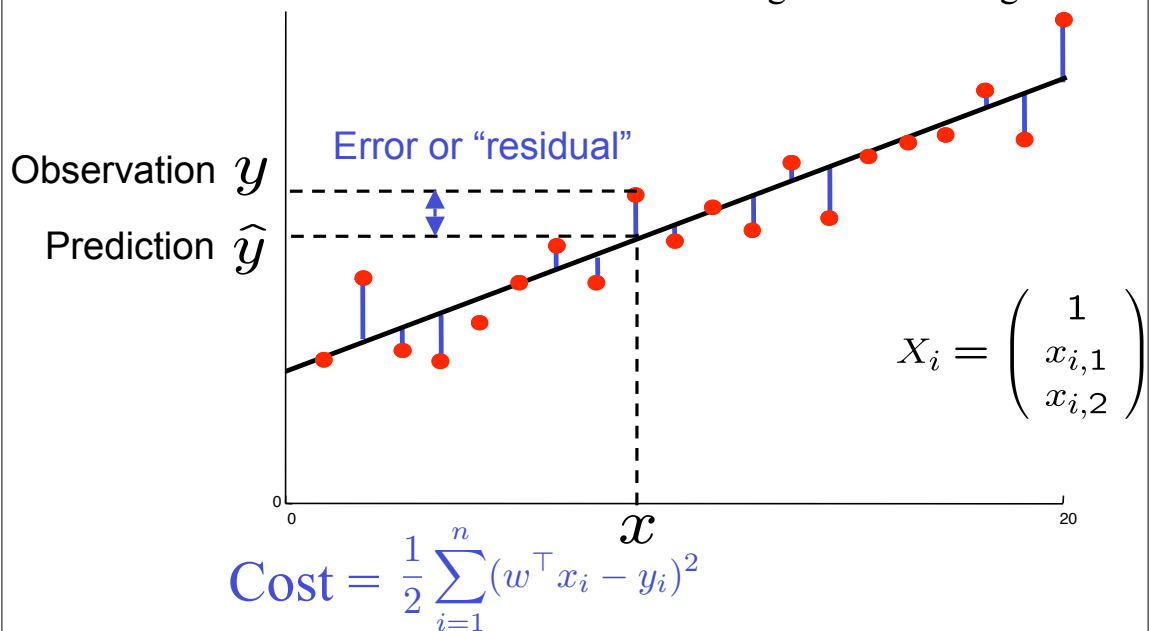
where w is a parameter to be estimated and we have used the standard convention of letting the first component of x be 1.



LMS Algorithm

(Least Mean Squares)

In order to clarify what we mean by a good choice of w , we will define a cost function for how well we are doing on the training data:



LMS Algorithm

(Least Mean Squares)

The best choice of w is the one that minimizes our cost function

$$E = \frac{1}{2} \sum_{i=1}^n (w^\top x_i - y_i)^2 = \sum_{i=1}^n E_i$$

In order to optimize this equation, we use standard gradient descent

$$w := w - \alpha \frac{\partial}{\partial w} E$$

where

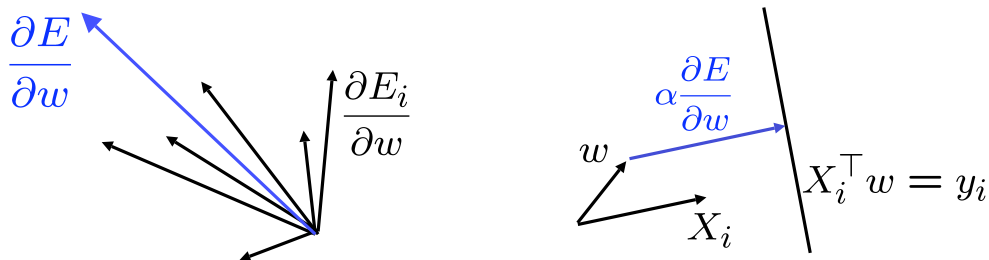
$$\frac{\partial}{\partial w} E = \sum_{i=1}^n \frac{\partial}{\partial w} E_i \quad \text{and} \quad \frac{\partial}{\partial w} E_i = \frac{1}{2} \frac{\partial}{\partial w} (w^\top x_i - y_i)^2 = (w^\top x_i - y_i) x_i$$

LMS Algorithm

(Least Mean Squares)

The LMS algorithm is an online method that performs the following update for each new data point

$$\begin{aligned} w^{t+1} &:= w^t - \alpha \frac{\partial}{\partial w} E_i \\ &= w^t + \alpha (y_i - x_i^\top w) x_i \end{aligned}$$



LMS, Logistic regression, and Perceptron updates

- LMS

$$w^{t+1} := w^t + \alpha(y_i - x_i^\top w)x_i$$

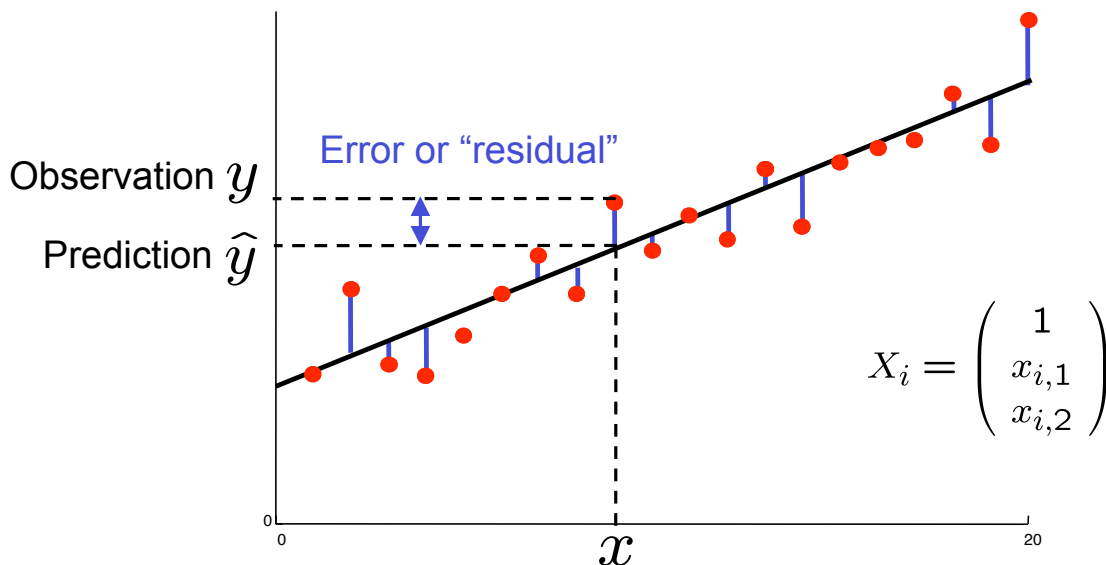
- Logistic Regression

$$w^{t+1} := w^t + \alpha(y_i - f_w(x_i))x_i$$

- Perceptron

$$w^{t+1} := w^t + \alpha(y_i - f_w(x_i))x_i$$

Ordinary Least Squares (OLS)



$$\text{Cost} = \frac{1}{2} \sum_{i=1}^n (w^\top x_i - y_i)^2$$

Minimize the sum squared error

$$\begin{aligned} E &= \frac{1}{2} \sum_{i=1}^n (w^\top x_i - y_i)^2 \\ &= \frac{1}{2} \|Xw - y\|_2^2 \\ &= \frac{1}{2} (w^\top X^\top Xw - 2y^\top Xw + y^\top y) \end{aligned}$$

$$X = \begin{pmatrix} -x_1^\top & - \\ -x_2^\top & - \\ \dots & \end{pmatrix} \begin{matrix} \updownarrow \\ n \end{matrix} \begin{matrix} \leftarrow \\ d \end{matrix}$$

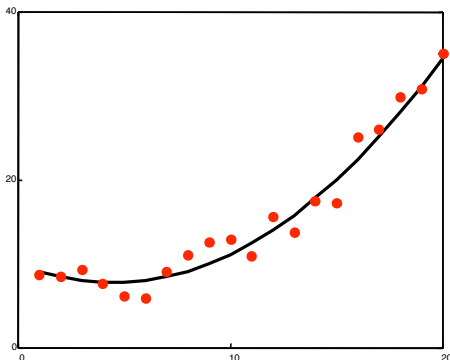
$$\frac{\partial}{\partial w} E = X^\top Xw - X^\top y$$

Setting the derivative equal to zero gives us the *Normal Equation*

$$\begin{aligned} X^\top Xw &= X^\top y \\ w &= (X^\top X)^{-1} X^\top y \end{aligned}$$

Beyond lines and planes

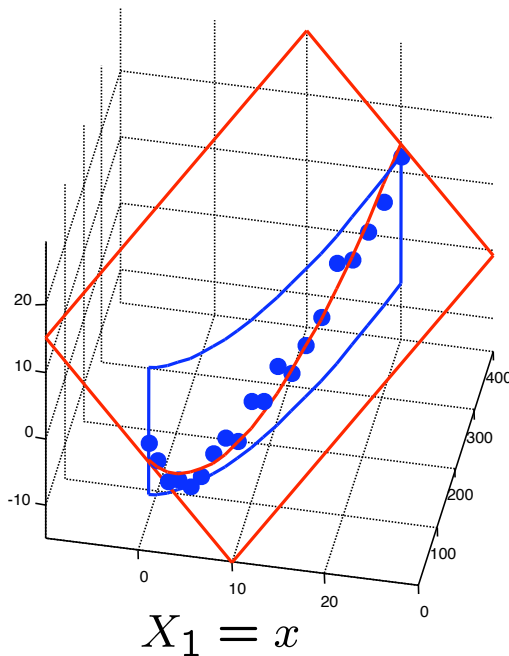
$$\hat{y}_i = w_0 + w_1 x_i + w_2 x_i^2$$



still linear in w

everything is the same with $X_i = \begin{pmatrix} 1 \\ x_i \\ x_i^2 \end{pmatrix}$

Geometric interpretation



$$\hat{y} = w_0 + w_1x + w_2x^2$$

$$X_2 = x^2$$

[Matlab demo]

Ordinary Least Squares [summary]

Given examples $(x_i, y_i)_{i=1\dots n}$

$$\text{Let } X_i^\top = (f_1(x_i) \quad f_2(x_i) \quad \dots \quad f_d(x_i))$$

$$\text{For example } X_i^\top = (1 \quad x_{i,1} \quad x_{i,2} \quad x_{i,1}^2 \quad x_{i,2}^2 \quad x_{i,1}x_{i,2})$$

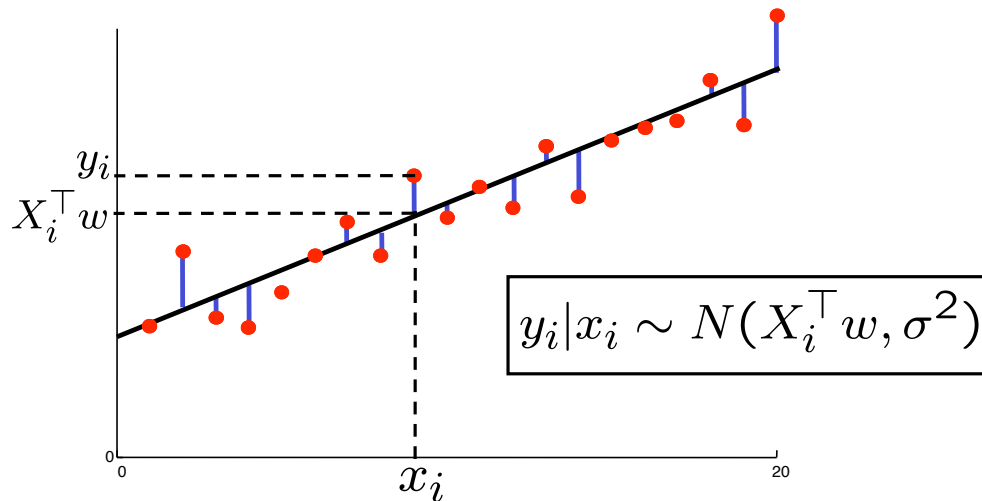
$$\text{Let } X = \begin{pmatrix} -X_1^\top - \\ -X_2^\top - \\ \dots \end{pmatrix} \begin{matrix} \updownarrow \\ n \\ \updownarrow \end{matrix} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \end{pmatrix}$$

$\leftarrow d \rightarrow$

Minimize $\|Xw - y\|_2^2$ by solving $(X^\top X)w = X^\top y$

Predict $\hat{y}_{n+1} = X_{n+1}^\top w$

Probabilistic interpretation



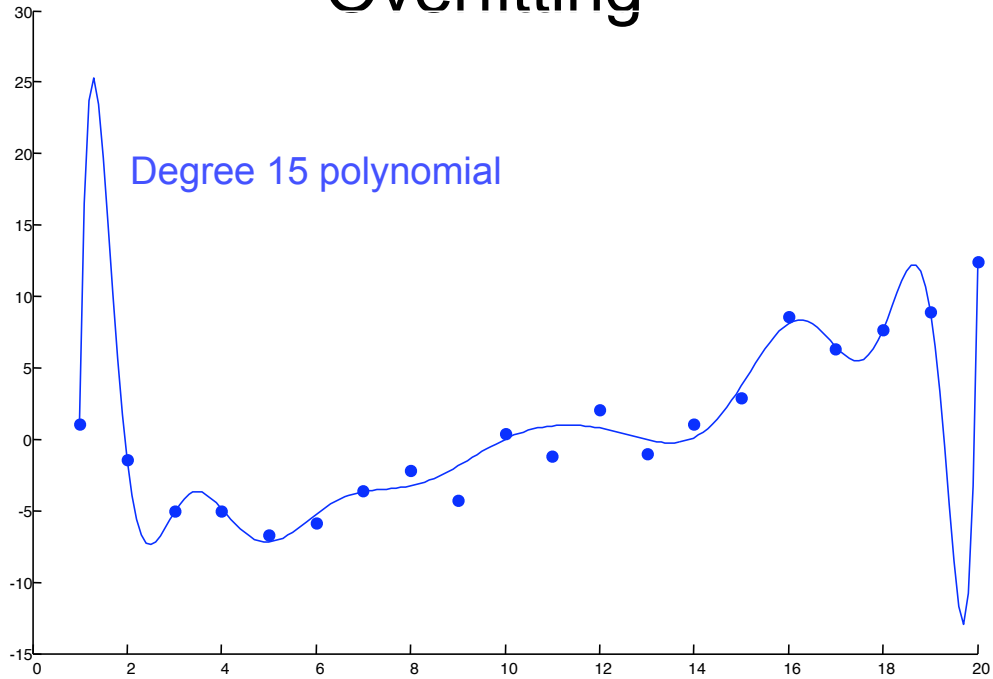
$$\text{Likelihood } L = \prod_i \exp -\frac{1}{2\sigma^2} (X_i^\top w - y_i)^2 = \exp -\frac{1}{2\sigma^2} \sum_i (X_i^\top w - y_i)^2$$

$$\operatorname{argmax}_w L = \operatorname{argmin}_w E$$

Outline

- Ordinary Least Squares Regression
 - Online version
 - Normal equations
 - Probabilistic interpretation
- **Overfitting and Regularization**
- Overview of additional topics
 - L_1 Regression
 - Quantile Regression
 - Kernel Regression and LWR
 - Spline Regression

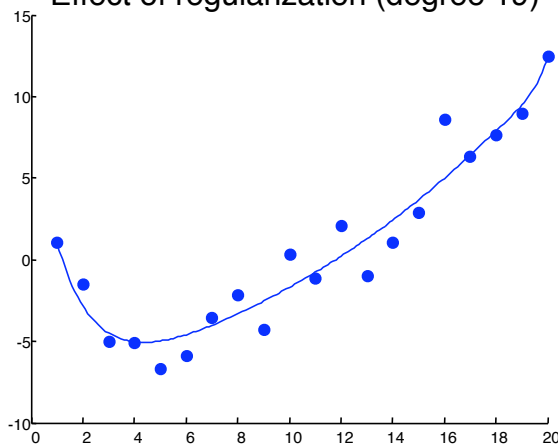
Overfitting



[Matlab demo]

Ridge Regression (Regularization)

Effect of regularization (degree 19)



$$A = X^T X$$

$$b = X^T y$$

with ϵ "small"

Minimize $\frac{1}{2} \|Xw - y\|_2^2 + \epsilon \|w\|_2^2$ by solving $(A + \epsilon I)w = b$

Probabilistic interpretation

Likelihood $y_i|x_i \sim N(X_i^\top w, \sigma^2)$

Prior $w \sim N\left(0, \frac{\sigma^2}{\epsilon}\right)$

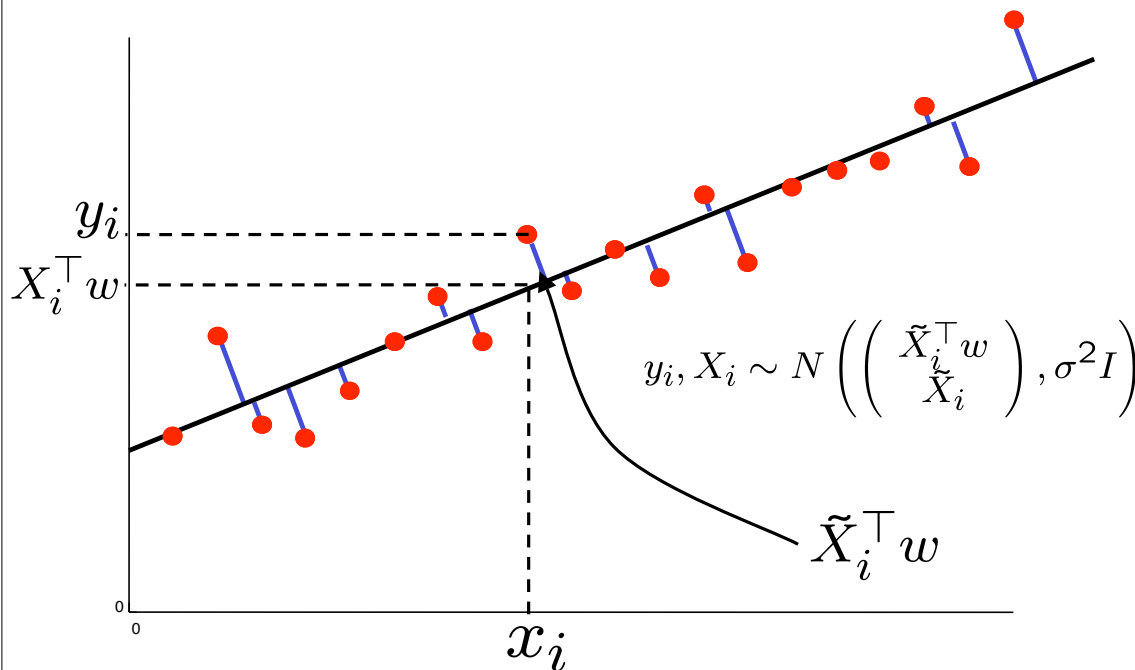
Posterior $P(w|x_1, \dots, x_n) = \frac{P(w, x_1, \dots, x_n)}{P(x_1, \dots, x_n)}$
 $\propto P(w, x_1, \dots, x_n)$

$$P(w, x_1, \dots, x_n) = \exp\left\{-\frac{\epsilon}{2\sigma^2}\|w\|_2^2\right\} \prod_i \exp\left\{-\frac{1}{2\sigma^2}(X_i^\top w - y_i)^2\right\}$$
$$= \exp\left\{-\frac{1}{2\sigma^2}\left[\epsilon\|w\|_2^2 + \sum_i (X_i^\top w - y_i)^2\right]\right\}$$

Outline

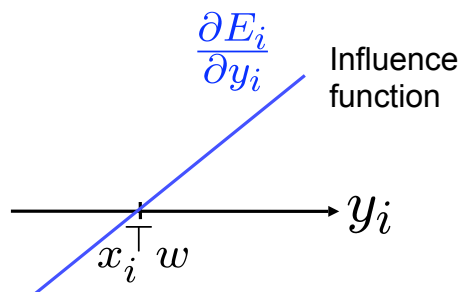
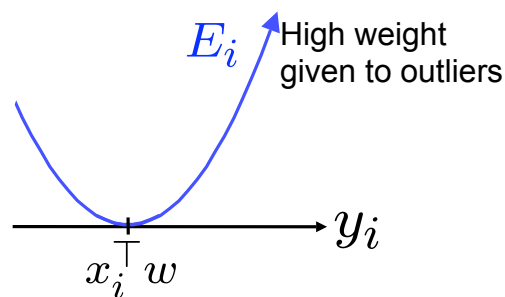
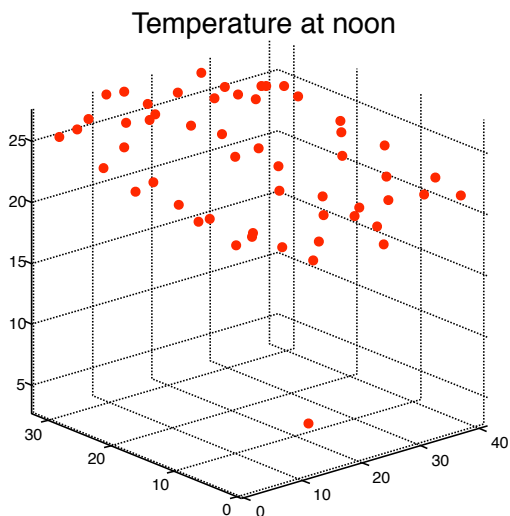
- Ordinary Least Squares Regression
 - Online version
 - Normal equations
 - Probabilistic interpretation
- Overfitting and Regularization
- Overview of additional topics
 - L_1 Regression
 - Quantile Regression
 - Kernel Regression and LWR
 - Spline Regression

Errors in Variables (Total Least Squares)



Sensitivity to outliers

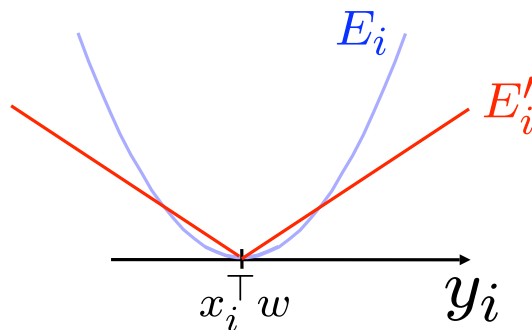
$$E = \sum_i (x_i^T w - y_i)^2 = \sum_i E_i$$



L₁ Regression

$$E' = \sum_i |x_i^\top w - y_i|$$

$$= \sum_i E'_i$$

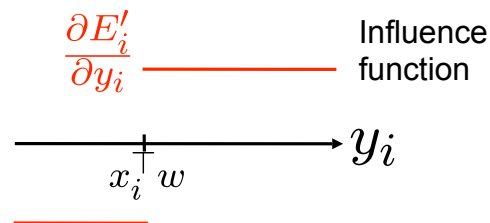


Linear program

$$\min_{w,c} \sum_i c_i$$

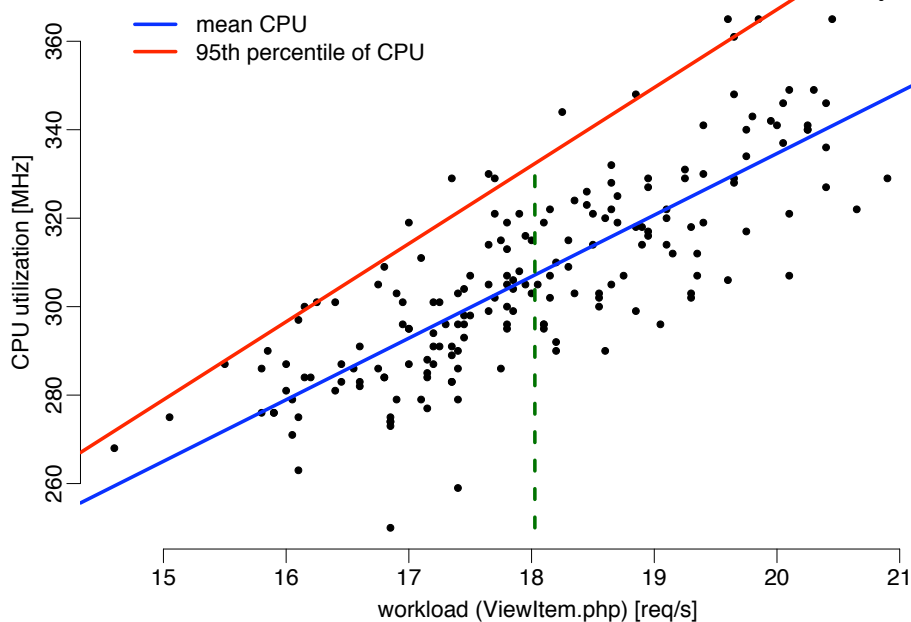
$$\text{s.t.} \quad x_i^\top w - y_i \leq c_i \quad \forall i$$

$$y_i - x_i^\top w \leq c_i \quad \forall i$$



[Matlab demo]

Quantile Regression



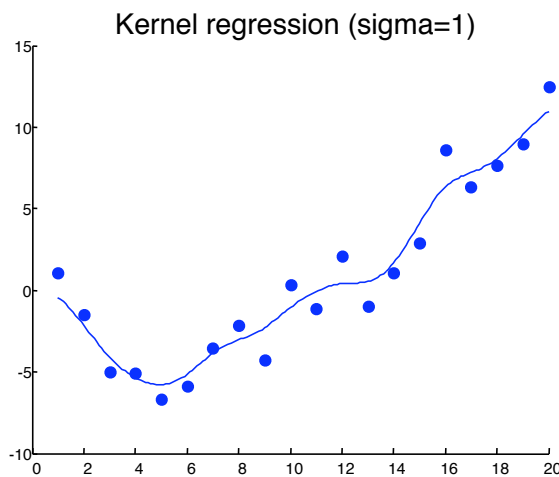
Slide courtesy of Peter Bodik

Kernel Regression and Locally Weighted Linear Regression

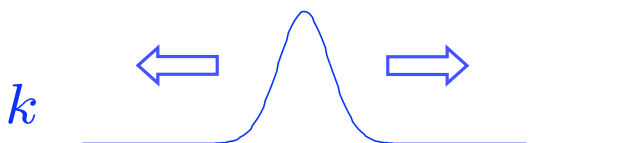
- **Kernel Regression:**
Take a very very conservative function approximator called AVERAGING. Locally weight it.
- **Locally Weighted Linear Regression:**
Take a conservative function approximator called LINEAR REGRESSION. Locally weight it.

Slide from Paul Viola 2003

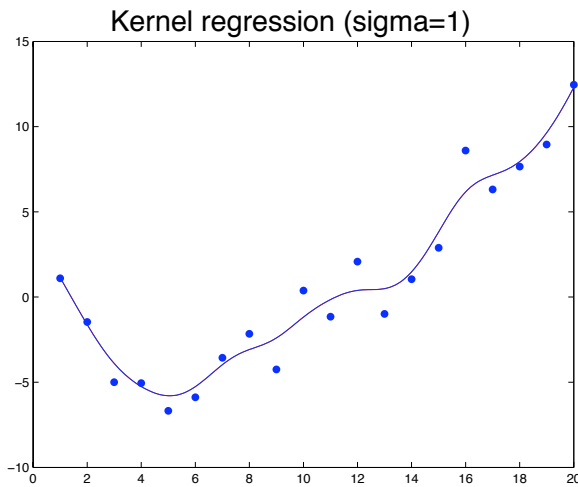
Kernel Regression



$$\hat{y}(x) = \frac{\sum_i y_i k(x_i - x)}{\sum_i k(x_i - x)}$$



Locally Weighted Linear Regression (LWR)

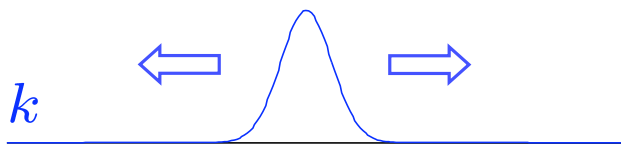


OLS cost function:

$$E = \frac{1}{2} \sum_{i=1}^n (w^\top x_i - y_i)^2$$

LWR cost function:

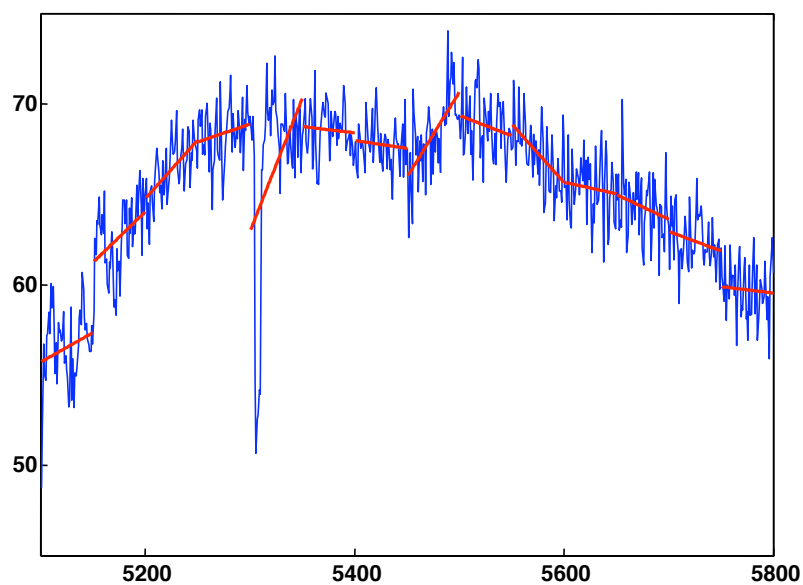
$$E' = \sum_{i=1}^n k(x_i - x) (w^\top x_i - y_i)^2$$



[Matlab demo]

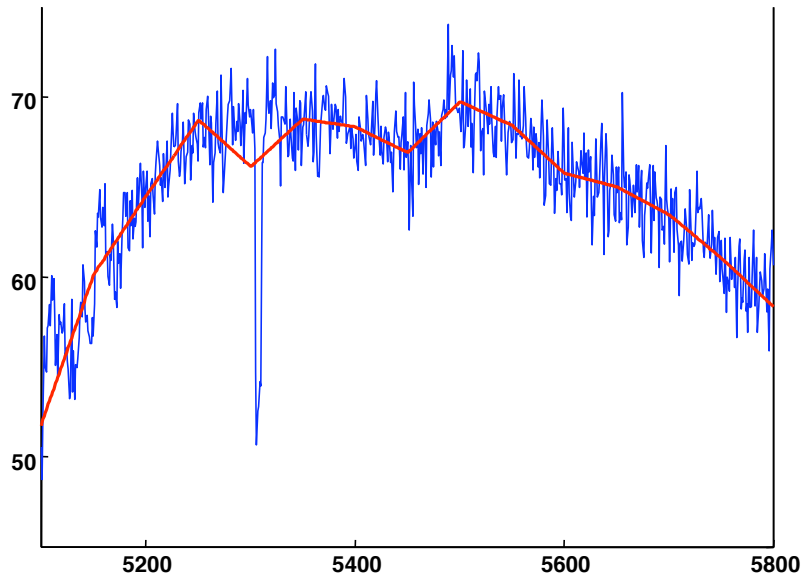
Spline Regression

Regression on each interval



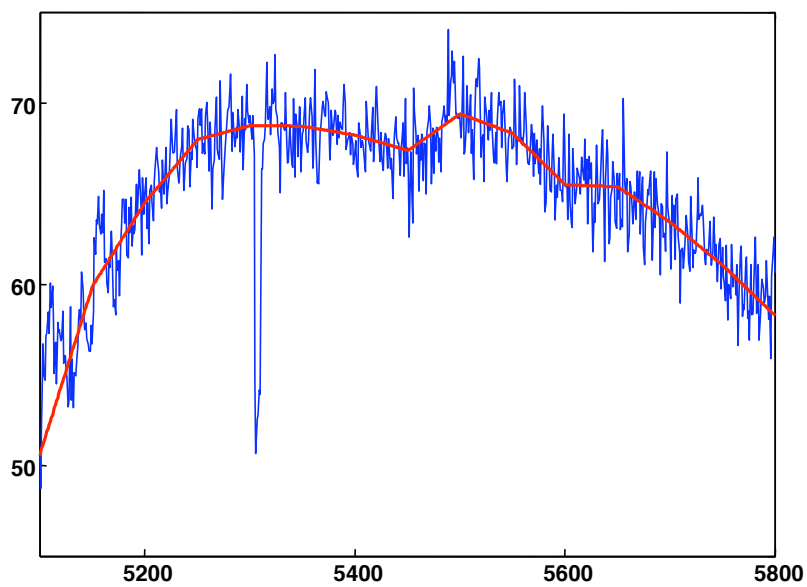
Spline Regression

With equality constraints



Spline Regression

With L_1 cost



Further topics

- Generalized Linear Models
- Gaussian process regression
- Feature Selection [after the break]